

# Trace analysis report SHARCNET

## General information

This is the trace analysis report (generated by reportgen.py) for the SharcNet system. The trace data was taken from the filename anon\_jobs.gwf, which contains job data obtained from. Below is a summary of the contents of the trace data:

- Date first entry: Wed Dec 21 02:55:33 2005
- CPU time consumed by jobs: 3782y 15d 7h 37m 57s
- Number of sites in the system: 10
- Number of CPUs in the trace: 6828
- Number of jobs in the trace: 1195242
- Number of users in the trace: 412
- Number of groups in the trace: 1

## System-wide characteristics

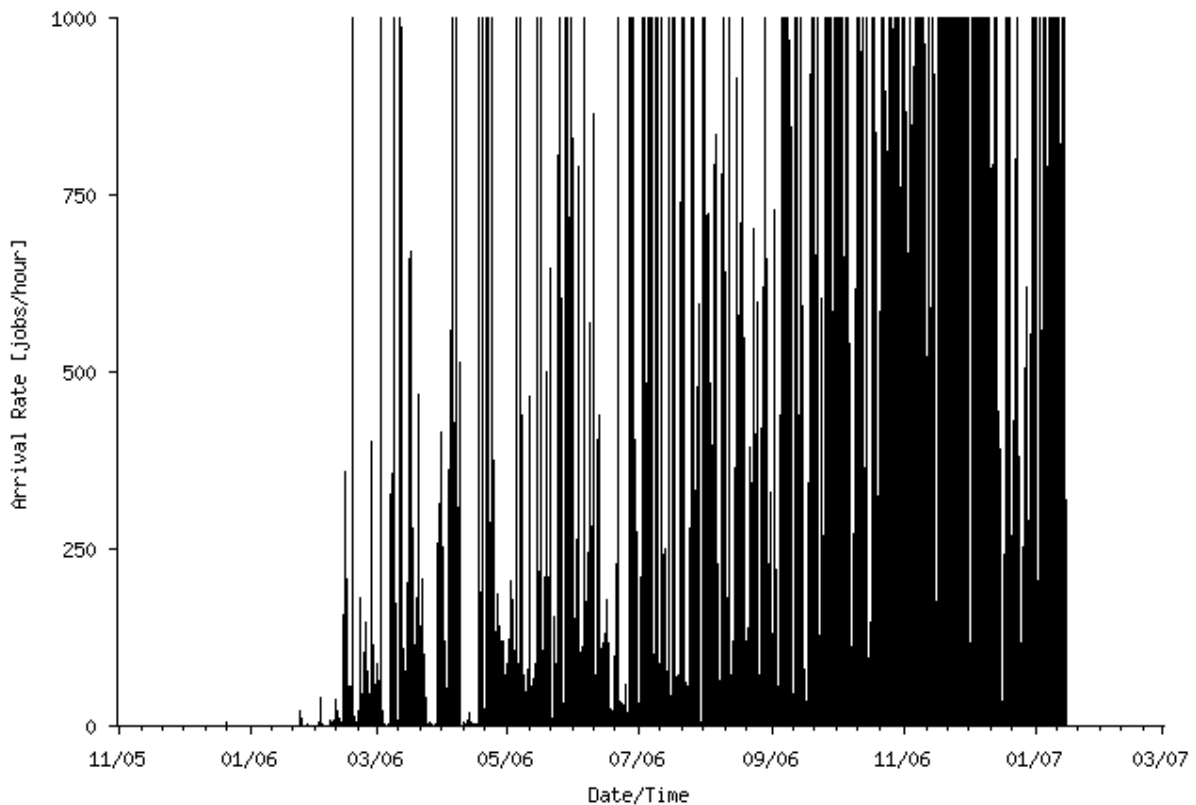
### System utilization

We define the overall system utilization as the ratio between the total CPU time consumed by users, and the total CPU time available to the users. We compute the total CPU time consumed by users as the sum of CPU time consumed by each job in the system; for failed jobs, only those that have effectively spent resource time are considered. We compute the total CPU time available as the number of CPUs multiplied by the duration of a fixed time interval, c.q. 10 minutes. Below we show the statistical properties of both the overall system utilization and the overall system for non-zero values, that is, excluding all intervals that have system utilization equal to zero. This excludes values that may account for downtime of the system. Unfortunately, utilization info for this trace is incomplete.

### Job arrival rate

We define the job arrival rate as the number of jobs that are submitted to the system in a fixed time interval. We compute the arrival rate for every hour by counting the all jobs that are recorded in the trace during that hour. This includes failed jobs and jobs that are cancelled before execution. Below we list the time periods in which the highest number of jobs were submitted to the system. We also summarize statistical properties for all job arrival rate values, and the statistical properties for arrival rate higher than zero. This excludes time periods that may account to downtime of the system. Figure 1 shows Overall job arrival rate during hourly intervals.

# SharcNet



**Figure 1 shows Overall job arrival rate during hourly intervals.**

Busiest time periods in terms of number of job submissions

- Busiest day: 2006-12-14
- Busiest week: 2006-44
- Busiest month: 2006-11

Overall job arrival metrics

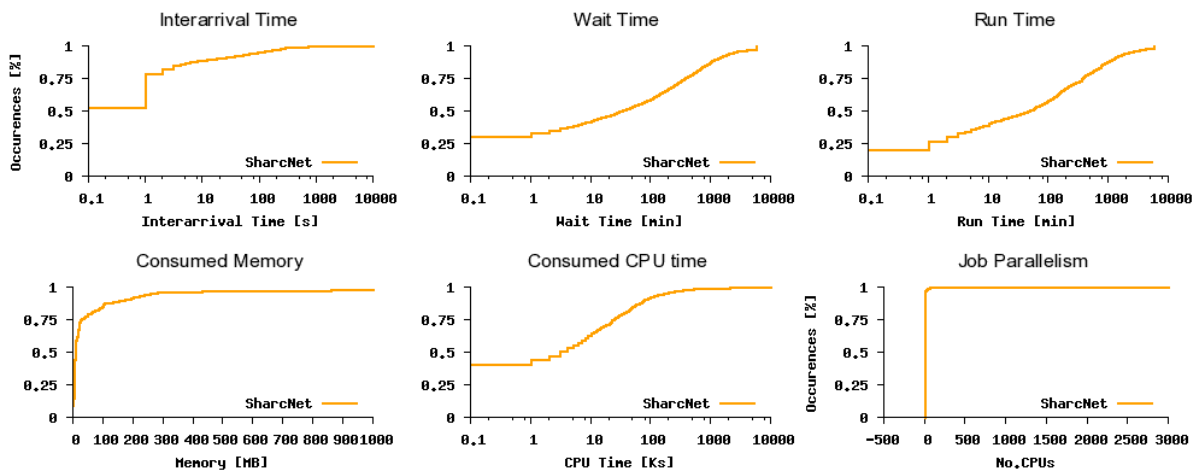
- Minimum: 0.00 jobs/hour
- Maximum: 22334.00 jobs/hour
- Average: 127.33 jobs/hour

Overall job arrival metrics for non-zero values

- Minimum: 2.00 jobs/hour
- Maximum: 22334.00 jobs/hour
- Average: 177.28 jobs/hour

## Job characteristics

We compute three important characteristics of jobs in the trace: number of CPUs used, the runtime of the job and the amount of memory used. Below we summarize the statistical properties for single jobs in the trace. We do not include jobs that were cancelled before execution, because those jobs did not consume resources from the system. Figure 2 shows CDFs of the most important job characteristics.



**Figure 2: CDFs of the most important job characteristics**

#### Number of CPUs used by a single job

- Minimum: 1 processors
- Maximum: 800 processors
- Average: 1.495 processors
- Standard deviation: 6.164
- Coefficient of variation: 4.122

#### Runtime of a single job

- Minimum: 0.00 seconds
- Maximum: 13908398.00 seconds
- Average: 31964.26 seconds
- Standard deviation: 117088.400
- Coefficient of variation: 3.663

#### Memory usage of a single job

- Minimum: 0.00 MB
- Maximum: 32021.50 MB
- Average: 81.28 MB
- Standard deviation: 466.176
- Coefficient of variation: 5.735

### Sequential vs. Parallel jobs

Below we summarize the resource usage of all sequential and all parallel jobs, that is all jobs that use more than one processor. First we calculate the number of sequential jobs and the number of parallel jobs that are submitted to the system. Furthermore, we compute the consumed CPU time by multiplying the runtime of a job by the number of processors allocated to the job. Again, this is divided into parallel and sequential jobs. For the number of jobs and the consumed CPU time, the percentage of all jobs is displayed

#### Number of jobs

- Sequential: 1064915 jobs (89.10 percent)
- Parallel: 117679 jobs (9.85 percent)

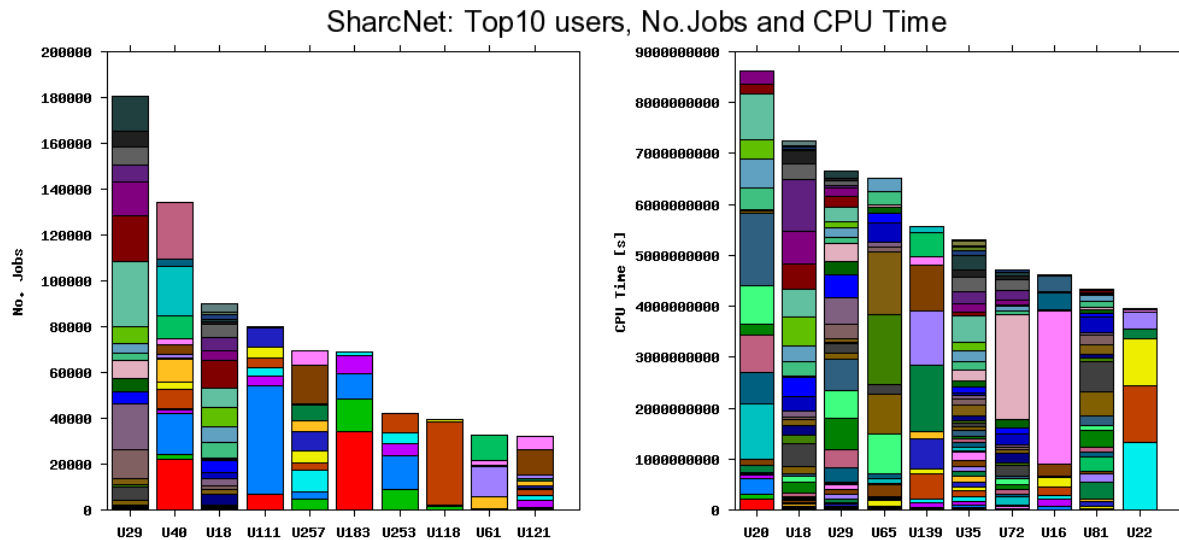
## Consumed CPU Time

- Sequential: 33840190929 seconds (28.37 percent)
- Parallel: 85629038048 seconds (71.79 percent)

## User and group characteristics

### User characteristics

Figure 3 shows The number of submitted jobs and the consumed CPU time by user.



**Figure 3: The number of submitted jobs (left) and consumed CPU time (right) by user. Only the top 10 users are displayed. The horizontal axis depicts the user's rank. The vertical axis shows the cumulated values, and the breakdown per week. Users have the same labels in the left and right sub-graphs**

### Top 10 users by number of job submitted to the system

Table 1 shows Top 10 users by number of jobs submitted to the system.

| Table 1 |        |                |            |
|---------|--------|----------------|------------|
| Rank    | UserID | Number of jobs | Percentage |
| 1       | U29    | 180456         | 15.10%     |
| 2       | U40    | 134321         | 11.24%     |
| 3       | U18    | 89777          | 7.51%      |
| 4       | U111   | 79843          | 6.68%      |
| 5       | U257   | 69657          | 5.83%      |
| 6       | U183   | 68774          | 5.75%      |
| 7       | U253   | 41965          | 3.51%      |
| 8       | U118   | 39565          | 3.31%      |
| 9       | U61    | 32750          | 2.74%      |
| 10      | U121   | 32051          | 2.68%      |

| Table 1 |        |                |            |
|---------|--------|----------------|------------|
| Rank    | UserID | Number of jobs | Percentage |
| 11      | Other  | 426083         | 35.65%     |
| 12      | Total  | 1195242        | 100.00%    |

#### Job arrival

- Minimum: 0.00 jobs/hour
- Maximum: 22235.00 jobs/hour
- Average: 96.06 jobs/hour

#### Job characteristics

##### Number of CPUs used by a single job

- Minimum: 1 processors
- Maximum: 800 processors
- Average: 1.495 processors
- Standard deviation: 6.164
- Coefficient of variation: 4.122

##### Runtime of a single job

- Minimum: 0.00 seconds
- Maximum: 2578047.00 seconds
- Average: 25917.65 seconds
- Standard deviation: 77395.107
- Coefficient of variation: 2.986

##### Memory usage of a single job

- Minimum: 0.00 MB
- Maximum: 8034.30 MB
- Average: 39.66 MB
- Standard deviation: 129.646
- Coefficient of variation: 3.269

#### Top 10 users by consumed CPU time

Table 2 shows Top 10 users by consumed CPU time (in seconds).

| Table 2 |        |             |            |
|---------|--------|-------------|------------|
| Rank    | UserID | CPU seconds | Percentage |
| 1       | U20    | 8628124335  | 7.23%      |
| 2       | U18    | 7245720464  | 6.08%      |
| 3       | U29    | 6662138572  | 5.59%      |
| 4       | U65    | 6506429650  | 5.46%      |
| 5       | U139   | 5554634120  | 4.66%      |
| 6       | U35    | 5295688828  | 4.44%      |
| 7       | U72    | 4704266918  | 3.94%      |
| 8       | U16    | 4611394976  | 3.87%      |
| 9       | U81    | 4338218963  | 3.64%      |
| 10      | U22    | 3953341652  | 3.31%      |

**Table 2**

| Rank | UserID | CPU seconds  | Percentage |
|------|--------|--------------|------------|
| 11   | Other  | 61770516999  | 51.79%     |
| 12   | Total  | 119270475477 | 100.00%    |

Job arrival

- Minimum: 0.00 jobs/hour
- Maximum: 5307.00 jobs/hour
- Average: 37.75 jobs/hour

Job characteristics

**Number of CPUs used by a single job**

- Minimum: 1 processors
- Maximum: 800 processors
- Average: 1.495 processors
- Standard deviation: 6.164
- Coefficient of variation: 4.122

**Runtime of a single job**

- Minimum: 0.00 seconds
- Maximum: 6717700.00 seconds
- Average: 49506.90 seconds
- Standard deviation: 154572.722
- Coefficient of variation: 3.122

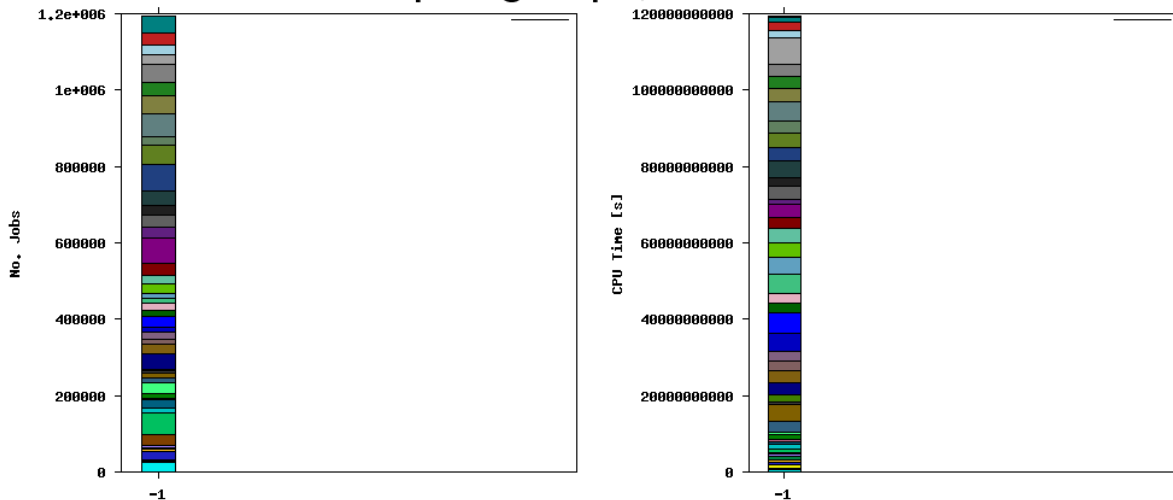
**Memory usage of a single job**

- Minimum: 0.00 MB
- Maximum: 29765.63 MB
- Average: 23.12 MB
- Standard deviation: 411.601
- Coefficient of variation: 17.802

Group characteristics

Figure 4 shows The number of submitted jobs and consumed CPU time by group.

**SharcNet: Top10 groups, No.Jobs and CPU Time**



**Figure 4: The number of submitted jobs (left) and consumed CPU time (right) by group. Only the top 10 groups are displayed. The horizontal axis depicts the groups rank. The vertical axis shows the cumulated values, and the breakdown per week. Groups have the same labels in the left and right sub-graphs**

Table 3 shows Top 10 groups by number of jobs submitted to the system.

| Table 3 |         |                |            |
|---------|---------|----------------|------------|
| Rank    | GroupID | Number of jobs | Percentage |
| 1       | 1       | 1195242        | 100.00%    |
| 2       | Other   | 0              | 0.00%      |
| 3       | Total   | 1195242        | 100.00%    |

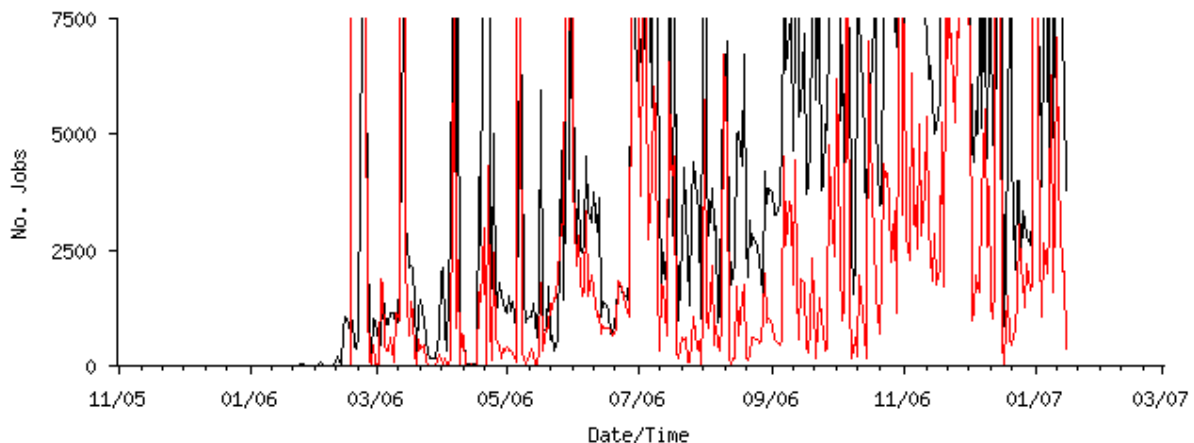
Table 4 shows Top 10 Groups by consumed CPU time (in seconds).

| Table 4 |         |              |            |
|---------|---------|--------------|------------|
| Rank    | GroupID | CPU seconds  | Percentage |
| 1       | 1       | 119270475477 | 100.00%    |
| 2       | Other   | 0            | 0.00%      |
| 3       | Total   | 119270475477 | 100.00%    |

## Performance analysis

### Waiting and running jobs

Figure 5 shows The number of running and of waiting jobs during hourly intervals. The vertical axis is limited to 7500 for better visibility.



**Figure 5: The number of running and of waiting jobs during hourly intervals. The vertical axis is limited to 7500 for better visibility**

We compute the number of running and waiting jobs by considering a fixed time interval. In each time interval, we count in the trace the amount of jobs that have been submitted but not yet started, that is, waiting. We also count the number of jobs that have been submitted, and have started executing in the time interval, but did not finish executing, and thus are running. Below we show the values for an interval value of 3600 seconds, summarized in amounts per day. Also the summary for values higher than zero are displayed, which excludes the possible effect of downtime of the system.

#### Number of waiting jobs per day

- Minimum: 0 jobs
- Maximum: 28056 jobs
- Average: 2507.78 jobs

#### Number of waiting jobs per day (non-zero values)

- Minimum: 1 jobs
- Maximum: 28056 jobs
- Average: 2952.11 jobs

#### Number of running jobs per day

- Minimum: 0 jobs
- Maximum: 34174 jobs
- Average: 4360.95 jobs

#### Number of running jobs per day (non-zero values)

- Minimum: 1 jobs
- Maximum: 34174 jobs
- Average: 4856.51 jobs

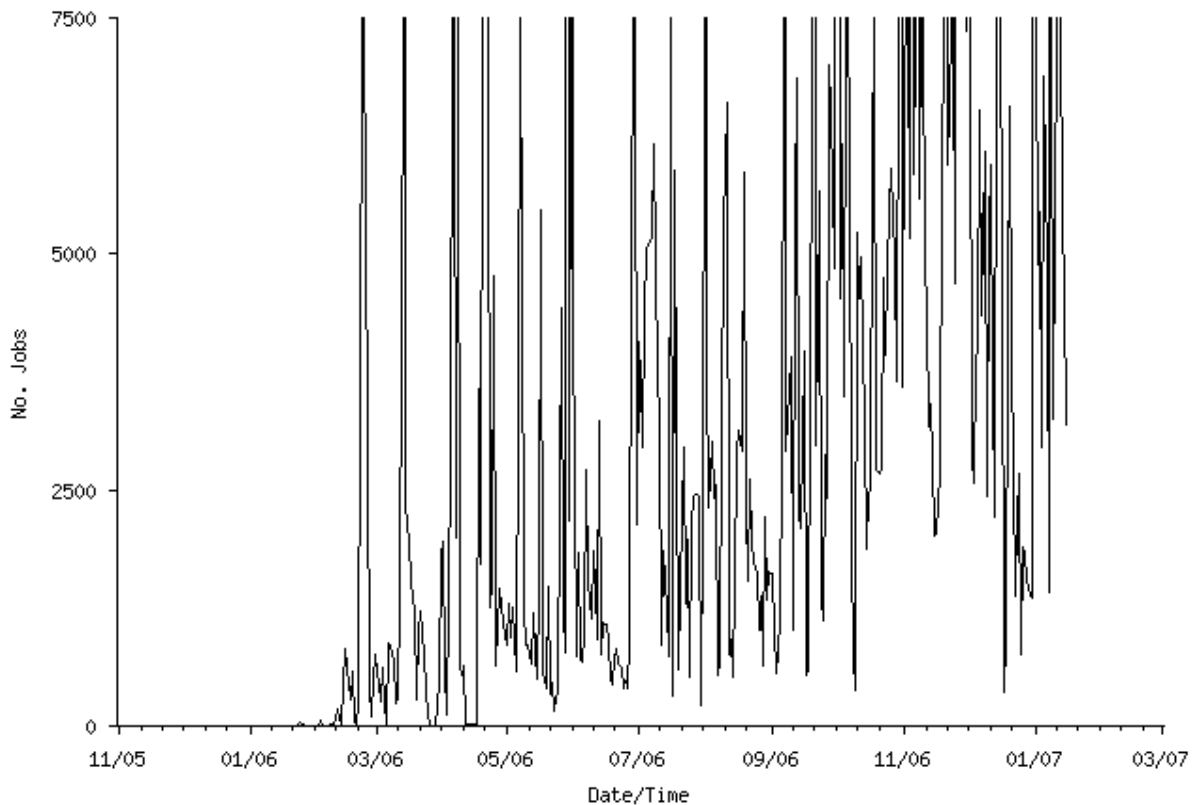
## Throughput

We compute the job throughput by considering a fixed time interval. In each time interval, we count in the trace the amount of jobs that have been submitted, started and finished executing. Below we show the values for an interval value of 3600 seconds, summarized in amounts per day. Also the summary for values higher than zero are displayed, which excludes the possible effect of downtime of the system.

Figure 6 shows Throughput during hourly intervals. The vertical axis of each individual site graph is limited to 7500 for better visibility.



# SharcNet



**Figure 6: Throughput during hourly intervals. The vertical axis of each individual site graph is limited to 7500 for better visibility**

## Throughput per day

- Minimum: 0 jobs
- Maximum: 31423 jobs
- Average: 3047.57 jobs

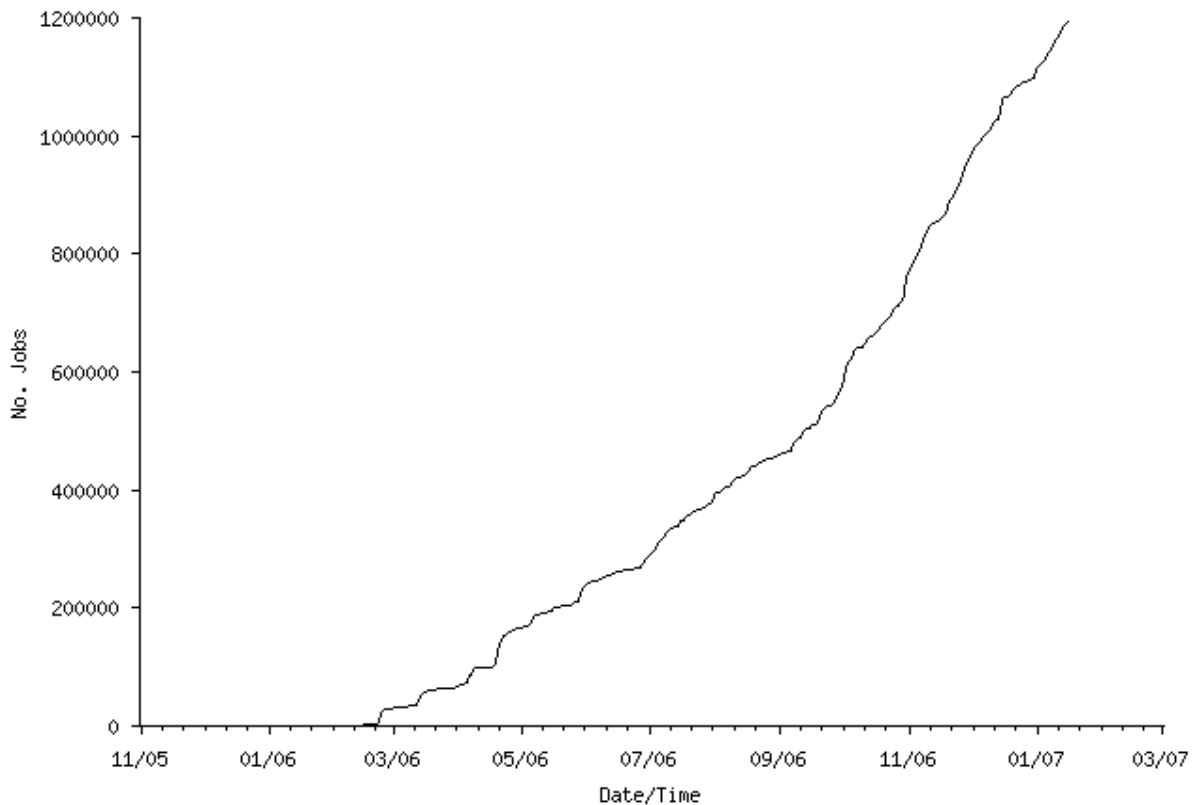
## Throughput per day (non-zero values)

- Minimum: 1 jobs
- Maximum: 31423 jobs
- Average: 3403.56 jobs

## Completed jobs

Figure 7 shows The number of completed jobs during hourly intervals.

# SharcNet



**Figure 7: The number of completed jobs during hourly intervals**

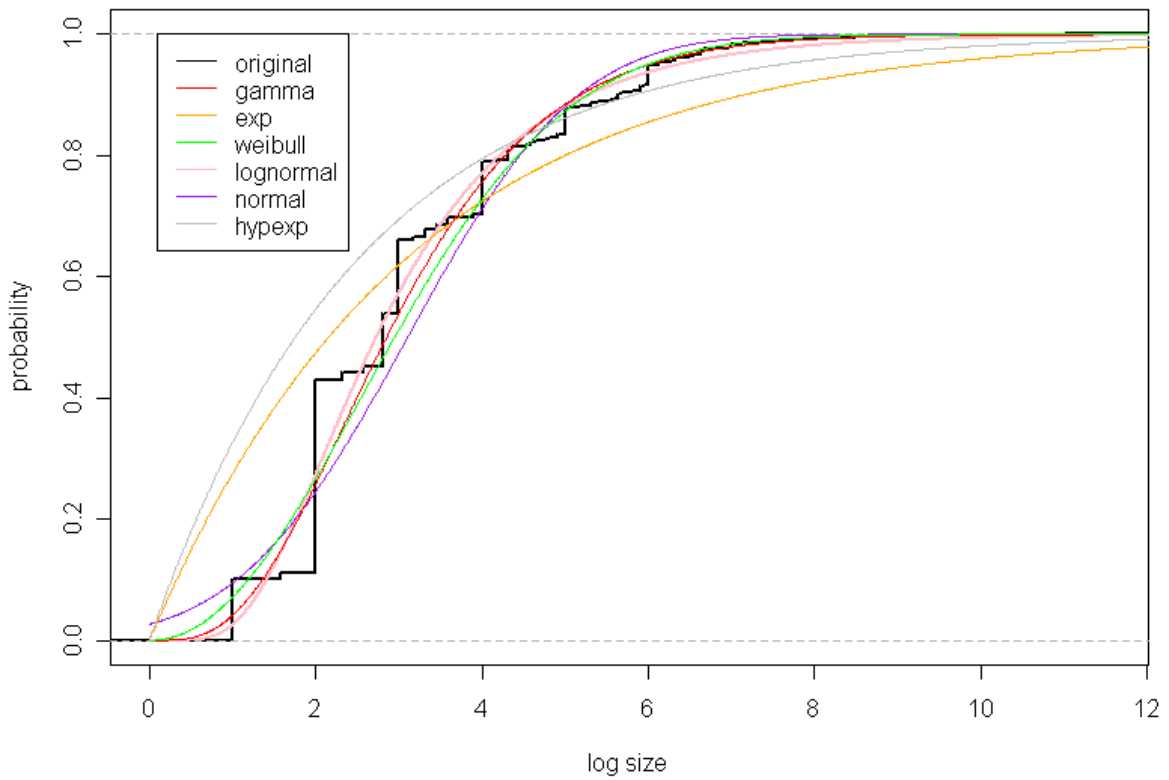
## Workload model

This section contains the workload model for the analyzed trace. The workload model consists of several parameters: job size, job runtime, requested runtime and interarrivals of jobs. These parameters are modeled by fitting well-known distributions to the data obtained from the trace. In all cases, first a logarithmic transformation was performed on the dataset to diminish the effect of outliers and speed up the modelling process. The fitting was performed using the maximum likelihood estimation method, which tries to maximize the log-likelihood function of each distribution given a dataset.

## Job size

Figure 8 shows Cumulative distribution function for the logarithm of the job sizes, with fitted distributions.

**Job size CDF with fitted distributions**

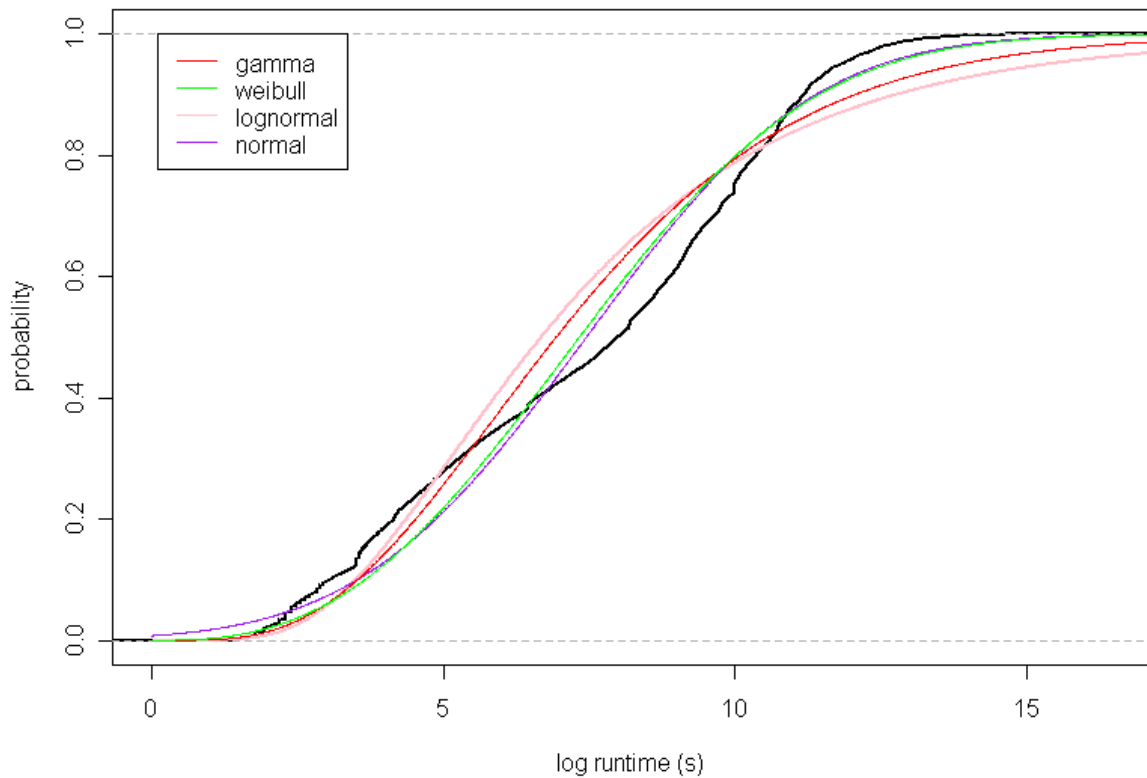


**Figure 8: Cumulative distribution function for the logarithm of the job sizes, with fitted distributions**

## Job runtime

Figure 9 shows Cumulative distribution function for the logarithm of the job runtimes, with fitted distributions.

### Job runtime CDF with fitted distributions



**Figure 9: Cumulative distribution function for the logarithm of the job runtimes, with fitted distributions**

### Parameters of fitted distributions

#### Lognormal distribution

- meanlog: 2.28868858386238
- sdlog: 0.476591812695073

#### Exponential distribution

- rate: 0.0919298782555215

#### Normal distribution

- sd: 4.24890758067712
- mean: 10.8778562419116

#### Weibull distribution

- shape: 2.7046976262594
- scale: 8.39041673471548

#### Hyperexponential distribution

- p: 6.23502620182541
- rate2: 0.106580708504889
- rate1: 0.220895993508368

## Gamma distribution

- shape: 4.72721961971207
- rate: 0.634731641752927

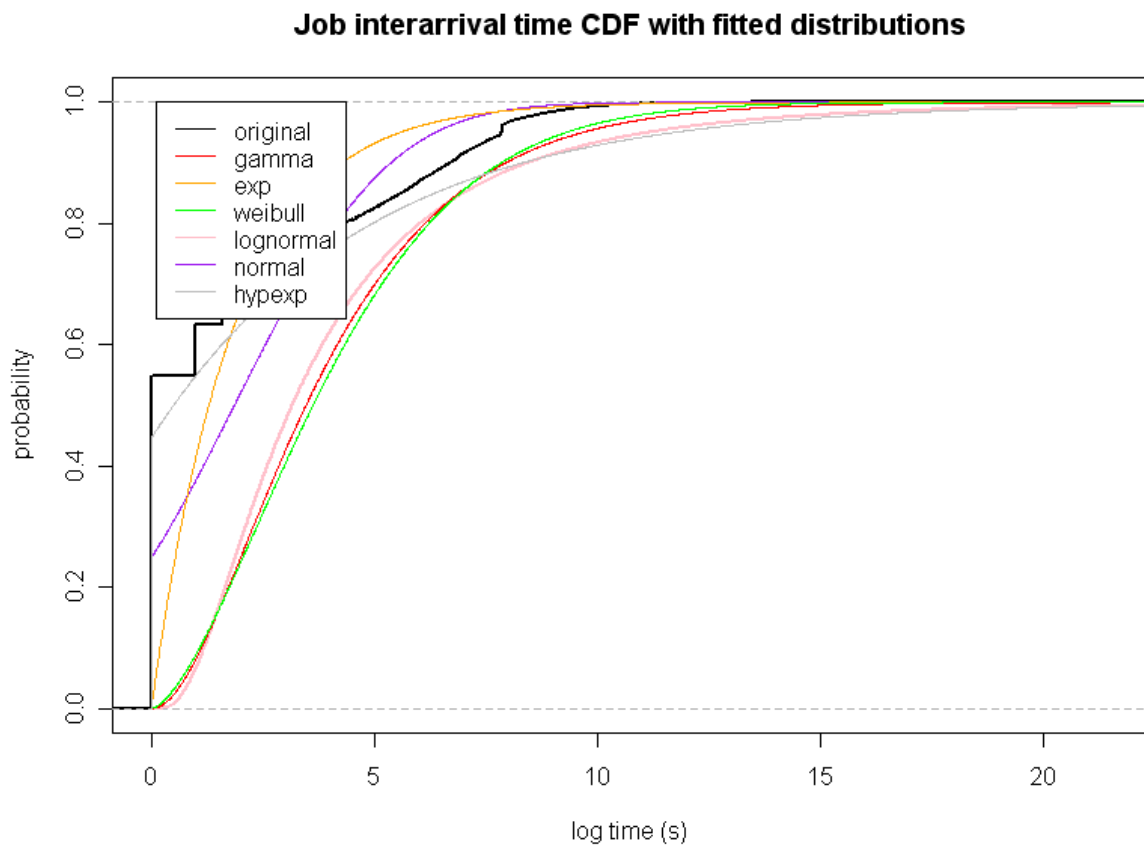
## Goodness-of-fit (Kolmogorov-Smirnov test)

Table 6 shows for each distribution the results for the Kolmogorov-Smirnov test, which gives a measure for the distance of the distribution to the original dataset (lower distance => better fit).

| Table 6          |                    |
|------------------|--------------------|
| Distribution     | Distance           |
| Lognormal        | 0.140903495377027  |
| Exponential      | 0.254038428047154  |
| Normal           | 0.081598730209193  |
| Weibull          | 0.0916160284600858 |
| Hyperexponential | 0.243453481913437  |
| Gamma            | 0.124966439992004  |

## Job interarrival

Figure 11 shows Cumulative distribution function for the logarithm of the job interarrival, with fitted distributions.



**Figure 11: Cumulative distribution function for the logarithm of the job interarrival, with fitted distributions**